

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

Document Summarizer For Word Processors

Inventor(s):

Ronald A. Fein

William B. Dolan

John Messerly

Edward J. Fries

Christopher A. Thorpe

Shawn J. Cokus

1 **TECHNICAL FIELD**

2 This invention relates to word processors, and more particularly, to
3 document summarizers for word processors.
4

5 **BACKGROUND OF THE INVENTION**

6 Many people are faced with the daunting task of reading large amounts of
7 electronic textual materials. In the computer age, people are inundated with
8 papers, memos, e-mail messages, reports, web pages, schedules, reference
9 materials, test results, and so on. Unfortunately, many documents do not begin
10 with summaries. Creation of summaries is tedious, requiring the author to re-read
11 the document, identify major themes, and distill the main points of the document
12 into a concise summary. Most authors never bother.

13 Summarizing a document is even more difficult and time-consuming for a
14 reader. The reader must first read the entire document (or at least skim it) to
15 understand the contents. The reader must then attempt to extract the document's
16 key points from unimportant details.

17 The problems associated with handling large volumes of un-summarized
18 documents are particularly acute for MIS (Management Information Systems)
19 personnel. These individuals are confronted daily with tasks of organizing,
20 managing, and retrieving documents from large databases. Imagine this typical
21 scenario. An MIS staff member receives a cryptic request to locate all documents
22 that pertain to a topic believed to have been discussed in a several company memos
23 written about three to four years ago. To accommodate this search request, the
24 MIS staff member must first perform a word search for the topic, and then
25 laboriously peruse each hit document in an effort to find the mysterious memos.

1 Without summaries, the staff member is forced to read large portions, if not all, of
2 each document before concluding whether the document is relevant or irrelevant.
3 Being forced to read unnecessary text leads to many wasted hours of the staff
4 member's time.

5 The problem is less critical, but still troubling, for individual users who are
6 browsing through the Internet or other networks to find documents on a related
7 topic. Upon locating a document, the user must either read the document online to
8 determine whether it is relevant (at the cost of additional online expenses), or
9 download the document for later review (at the risk of retrieving an irrelevant
10 document).

11 To help address these problems, computer-implemented document
12 summarizers have been developed to automatically summarize text-based
13 documents for the readers. The document summarizers examine an existing
14 document, and attempt to create an abstract or summary from the existing text.

15 Early development on document summarizers centered on statistical
16 approaches to creating summaries. One statistical approach is described in an
17 article by H.P. Luhn, entitled "The Automatic Creation of Literature Abstracts,"
18 which was published April 1958 in the IBM Journal at pages 159-165. The Luhn
19 technique assigns to each sentence a "significance" factor derived from an analysis
20 of its words. This factor is computed by ascertaining a cluster of words within a
21 sentence, counting the number of significant words contained in the cluster, and
22 dividing the square of this number by the total number of words in the cluster. The
23 sentences are then ranked according to their significance factor, with one or
24 several of the highest ranking sentences being selected to form the abstract.

1 Most, if not all, of the document summarizers in use today appear to employ
2 the Luhn technique. Examples of such summarizers include a Text Summariser
3 from BT (formerly British Telecom), Visual Recall from Xsoft Corporation (a
4 subsidiary of Xerox), and InText from Island Software.

5 Another approach to summarizing documents is described in an article by
6 Kenji Ono, et al., entitled "Abstract Generation Based on Rhetorical Structure
7 Extraction," which was published in Proceedings of the 15th International
8 Conference on Computational Linguistics, Vol. 1, at pages 344-348, for a
9 conference held Aug 5-9, 1994 in Kyoto, Japan. Their approach involved a
10 linguistic analysis which constructed rhetorical structures representing relations
11 between various chunks of sentences in the body of the section. The rhetorical
12 structure is represented by two levels: intra-paragraph, which analyzes the text
13 according to sentence units, and inter-paragraph, which analyzes the text using
14 paragraph units. Extraction of the rhetorical structure is accomplished using a
15 detailed and sophisticated five-step procedure. The Ono technique is unnecessarily
16 complicated for many situations where a rudimentary summary is all that is
17 desired.

18 In addition, this technique is highly genre-dependent, producing good
19 summaries only when the text is rich in superficial markers of its discourse
20 structure. It thus works relatively well on the academic prose examined by Ono et
21 al., but will fail on documents written in less formal prose.

22 When the summaries are created, conventional document summarizers
23 present the results to the reader in one of two formats. The first format is to
24 underline or otherwise highlight the sentences that are deemed to be part of the
25

1 summary. The second format is to show only the abstracted sentences in paragraph
2 or bullet format, without the accompanying text of the document.

3 One common problem with the conventional document summarizers is that
4 they are *reader*-based. These summarizers do not consider summary creation and
5 presentation from the perspective of the *author*.

6 Accordingly, there remains a need to provide an *author*-oriented
7 summarizer for a word processor that helps authors automatically create
8 summaries for their writings, and one which will produce a summary for any text
9 which is presented to it.

10 11 SUMMARY OF THE INVENTION

12 This invention concerns a document summarizer which is particularly
13 helpful in assisting authors in preparing summaries for documents, as well aiding
14 readers in their review of un-summarized documents. For a given text, the
15 document summarizer first performs a statistical analysis to generate a list of
16 ranked sentences for consideration in the summary. The summarizer counts how
17 frequently content words appear in a document and produces a table correlating the
18 content words with their corresponding frequency counts. A sentence score for
19 each sentence is derived by summing the frequency counts of the content words in
20 the sentence and dividing that sum by the number of the content words in the
21 sentence. The sentences are then ranked in order of sentence scores, with higher
22 ranking sentences having comparatively higher sentence scores and lower ranking
23 sentences having comparatively lower sentence scores.

24 Concurrent with the statistical analysis in the same pass through the
25 document, the document summarizer performs a cue-phrase analysis by consulting

1 a pre-compiled list of words and phrases which serve either as indicators of
2 discourse relationships between adjacent sentences in a document or as an
3 indicator of the overall importance of a particular sentence in a document. The
4 cue-phrase analysis compares the sentence string to this pre-compiled list of cue
5 phrases. Associated with each cue phrase are conditions which are used to
6 determine whether a sentence containing that cue phrase will be used in a
7 summary.

8 For instance, the list might contain words and phrases which depend on the
9 surrounding context of the document to properly understand the sentence. A
10 sentence that begins, "That is why..." or "In contrast to this..." depends on
11 statements made in the preceding sentence(s). The summarizer establishes a
12 condition that a sentence containing a dependent word or phrase may only be
13 included in the summary if the neighboring context from which the word or phrase
14 depends is also included in the summary.

15 The pre-compiled list also contains cue phrases whose presence in a
16 sentence will result in that sentence being excluded from the summary, no matter
17 how large its statistically-derived score might be. For instance, a sentence which
18 contains the phrase "as shown in Fig...." should not be included in a summary
19 because the referenced figure will not be present.

20 Following the statistical and cue-phrase analysis phases, the summarizer
21 creates a summary containing the higher ranked sentences. The summary may also
22 include a conditioned sentence (such as one that contains a dependent word or
23 phrase) if the conditions established for inclusion of the sentence have been
24 satisfied. However, the summary never includes prohibited sentences.
25

1 The summarizer inserts the sentence at the beginning of the document
2 before the start of the text, or in a new document, based on the user's choice. This
3 placement is convenient and useful to the author. The author is then free to revise
4 the summary as he/she wishes.

6 **BRIEF DESCRIPTION OF THE DRAWINGS**

7 Fig. 1 is a diagrammatic illustration of a computer loaded with a word
8 processing program having a document summarizer.

9 Fig. 2 is a flow diagram of steps in a computer-implemented method for
10 summarizing documents.

11 Figs. 3a and 3b show documents with summaries inserted therein to
12 illustrate two different display presentations of a summary.

14 **DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT**

15 Fig. 1 shows a computer 20 having a central processing unit (CPU) 22, a
16 monitor or display 24, a keyboard 26, and a mouse 28. Other input devices--such
17 as a track ball, joystick, and the like--may be substituted for or used in conjunction
18 with the keyboard and mouse. The CPU 22 is of standard construction, including
19 memory (disk, RAM, graphics) and a processor.

20 The computer 20 runs an operating system which supports multiple
21 applications. The operating system is stored in memory in the CPU 22 and
22 executes on the processor. The operating system is preferably a multitasking
23 operating system which allows simultaneous execution of multiple applications.
24 One example operating system is a Windows® brand operating system sold by
25 Microsoft Corporation, such as Windows® 95 or Windows NT™ or other

1 derivative versions of Windows®. However, other operating systems may be
2 employed, such as Mac™OS operating systems employed in Macintosh computers
3 manufactured by Apple Computer, Inc.

4 This invention concerns a document summarizer that can be implemented in
5 a word processing system. In the illustrated system, the word processing system is
6 implemented as a software application which is stored in the CPU memory or other
7 loadable storage medium and runs on the operating system of computer 20. One
8 example word processing application is Microsoft® Word from Microsoft
9 Corporation, which is modified with the document summarizer described herein.

10 It is noted that the word processing system might be implemented in other
11 ways. For instance, the word processing system might comprise a dedicated
12 typewriter machine with limited memory and processing capabilities (in
13 comparison to a personal computer) that is used almost exclusively for word
14 processing tasks. It is further noted that the document summarizer described
15 herein can be implemented in other programs, such as an Internet Web browser
16 (e.g., Internet Explorer from Microsoft Corporation), an e-mail program (e.g.,
17 WordMail and Exchange from Microsoft Corporation), and the like. However, for
18 discussion purposes, the document summarizer is described in the context of a
19 computer word processing program, such as Microsoft® Word.

20 When an author wishes to summarize a document, the author initiates the
21 document summarizer function on the word processing program. As used herein,
22 the term "document" means any image that contains text in a format intended for a
23 viewer or other computer program which will then present the text as intelligible
24 language. Examples of documents include conventional word processing
25 documents, e-mail messages, memoranda, web pages, and the like. The document

1 summarizer is activated through a pull down menu or soft button on the graphical
2 user interface window presented by the word processor. Upon activation, the
3 document summarizer begins processing a document to produce a summary.

4 Fig. 2 shows the general steps in a computer-implemented method for
5 summarizing a document that are carried out by the computer. The method is
6 described with additional reference to an example document containing a four-
7 sentence paragraph, which is summarized into a two-sentence summary. The
8 paragraph is given as follows:

9
10 The Internet is a great place to shop for
11 a computer. Manufacturers have web sites
12 describing their computers. One computer
13 manufacturer offers a money back guarantee.
14 That is why that manufacturer has so many
15 visits to its Internet web site.
16

17 In general, the document summarizing process involves three phases: a
18 statistical phase, a cue-phrase phase, and a presentation phase. The statistical and
19 cue-phrase phases are preferably conducted concurrently during a single pass
20 through a document. However, they can be performed sequentially as well, in any
21 order. In the statistical phase, the document summarizer begins by reading each
22 word and counting how frequently content words appear in a document (step 40 in
23 Fig. 2). "Content words" are those words which provide non-grammatical
24 meaning to a text. Nouns are good examples of content words. In the above
25

1 paragraph, content words include "Internet," "manufacturer," "computer," and so
2 forth.

3 Within the context of the summarizer, content words can be technically
4 defined as words that are not "stop words." In this context, the set of stop words
5 includes both grammatical function words (e.g. conjunctions, articles,
6 prepositions) and certain high frequency verbs and nouns (e.g. "get", "have")
7 which appear to contribute relatively little semantic content to a sentence. The
8 fundamental attribute of a stop word is that it does not directly contribute to the
9 theme of the document, and the document is extremely unlikely to be about the
10 stop word; therefore it should not be counted. The stop words are preferably
11 maintained in a list stored in memory. In this manner, the processor reads every
12 word, but only counts those words that do not appear on the stop word list. In the
13 above sample paragraph, the first sentence contains the stop words "The," "is,"
14 "a," "great," "to," "for," and "a."

15 During the pass through the document, the document summarizer checks for
16 morphological variants of the content words and converts them to their root form
17 (step 42). For example, the words "walking," "walked," and "walks" are all
18 morphological variants of the root form "walk." In this way, the root form and
19 associated variants are all counted as the same word. In the above example
20 paragraph, the words "computer" and "computers" are counted as the same word,
21 as are the words "manufacturer" and "manufacturers."

22 The summarizer also analyzes the words for possible phrase compression
23 (step 44). Sets of content words that appear repeatedly in the same order are
24 counted as if they are a single content word. For example, the word pair,
25 "Microsoft Corporation," if occurring a sufficient number of times in that exact

1 order, might be counted as a single word. The words in such phrases, if taken
2 separately, do not by themselves add any meaning to the sentence. Without phrase
3 compression, the words "Microsoft" and "Corporation" would each be counted
4 independently, a result which might undesirably skew the importance of the
5 sentences that contain them. In the above example paragraph, the phrase "web
6 site" occurs the same way on two occasions and might therefore be a candidate for
7 phrase compression. Also assume that the phrase "money back guarantee" is
8 compressed into one word phrase that is counted singly.

9 When all of the content words in the document are counted, the document
10 summarizer produces a table which correlates the content words with their
11 corresponding frequency counts (step 46). The content words can be ordered with
12 the most frequently occurring words appearing at the top of the table. Table 1
13 shows a ranking of content words from the above example document:

14
15 **Table 1: Rank of Content Words**

<u>Content Word</u>	<u>Frequency Count</u>
Computer	3
Manufacturer	3
Internet	2
web site	2
Place	1
Shop	1
money back guarantee	1
Visit	1

At step 48, the document summarizer derives a sentence score for individual sentences within the document according to their respective content words. Sentences with more content words that appear more frequently in the document are ranked higher than both sentences with fewer high-frequency content words and sentences with content words that appear less frequently in the document. More specifically, the document summarizer ranks the sentences according to their average word score. This value is derived by summing the frequency counts for all content words that appear in the sentence and dividing that tally by the number of the content words in the sentence. The sentence score is represented as follows:

$$\text{Sentence Score} = \text{Sum of Word Frequency Counts} \div \text{Number of Words}$$

The sentences are then ranked in order of their sentence scores (step 50 in Fig. 2). Higher ranking sentences have comparatively higher sentence scores and lower ranking sentences have comparatively lower sentence scores. Using the word counts in Table 1, the score for the first sentence in the example paragraph is 1.75, as follows:

$$\text{Sentence \#1} = [\text{Internet}(2) + \text{Place}(1) + \text{Shop}(1) + \text{Computer}(3)] \div 4 \text{ Words} = 1.75$$

Scores for the remaining three sentences are also computed. Table 2 shows the ranking for the four sentences in the example paragraph.

Table 2: Rank of Sentences

<u>Sentences</u>	<u>Score</u>
#2 Manufacturers have web sites describing their...	2.67
#3 One computer manufacturer offers a money back...	2.33
#4 That is why that manufacturer has so many visits to ...	2.00
#1 The Internet is a great place to shop for a computer.	1.75

It is noted that other techniques could be used to derive a sentence score. For example, the score might be calculated by dividing the total frequency count by the total number of all words (including stop words) in sentence. An alternative approach is to simply sum the content word counts, without any averaging. Additionally, arithmetic and statistical tricks can be used, such as basing the sentence score on a median score of a content word.

Steps 40-50 constitute the statistical phase of the summarizing method. Concurrent with the statistical phase, the document summarizer performs during the same pass through the document a cue-phrase analysis to exploit any explicit discourse markers present in the text. In general, the cue-phrase analysis seeks to identify phrases that might potentially render a sentence confusing or difficult to understand if included in the summary. In this implementation, the document summarizer compares the sentence string to a pre-compiled list of words and phrases (step 52).

Upon identification of words or phrases that appear on the list, the document summarizer designates the entire sentence as either "prohibited" or "conditioned." If a sentence is "prohibited," the document summarizer takes action to prevent the sentence from being included in the summary, regardless of

1 its sentence score (steps 54 and 56). If a sentence is deemed “conditioned,” the
2 document summarizer will only include the sentence in the summary if the
3 condition is met (steps 58 and 60). One example of a conditioned sentence is one
4 that depends on the previous sentence or surrounding context to understand its
5 meaning. A sentence that begins “He said...” is only clear if the reader knows
6 who “He” is. Accordingly, this sentence depends on a previous context and will
7 be used in the summary only if the previous sentence identifying “He” is also used
8 in the summary.

9 Table 3 shows example words and phrases from the pre-compiled cue-
10 phrase list that render a sentence as “prohibited” or “conditioned.”

11
12 **Table 3: Cue-Phrase List**

13 **Conditional Words or Phrases**

14 Sentence-initial Personal Pronouns: He, She, It, They, Their

15 Sentence-initial Demonstrative Pronouns: These, That, This, Those

16 Sentence-initial Quantifiers: All, Most, Many

17 Both, Which

18 Conjunction (i.e., And, Nor, But, Or, Yet, So, For)

19 Specific Reference (i.e., Such, That is)

20 Extension (i.e., Related to this)

21 Causation (i.e., Therefore, Thus, And so)

22 Contrast (i.e., However, Nonetheless, In spite of this)

23 Reinforcement (i.e., Indeed, Accordingly)

24 Supplementation (i.e., At any rate, In reply)

Prohibited Words or Phrases

Reference (i.e., In Fig. 1..., as shown in Chart A)

Applying the cue phrase analysis to the sample paragraph reveals that the fourth sentence is conditional because it contains the phrase "That is why..." This phrase is listed on the cue-phrase list as a depends-on-previous phrase, meaning that the phrase relies on a previous sentence for context. In this case, the preceding third sentence explains that one manufacturer offers a money back guarantee which is the supporting reason why the manufacturer is said, in the fourth sentence, to have many visits to its web sit. Were the fourth sentence to appear in a summary without the third sentence, a reader would not understand why the manufacturer has so many visits to its web site. Accordingly, the document summarizer sets a condition that the fourth sentence is only used in the summary if the third sentence is also used.

In this example, it turns out that even without the cue phrase list, the fourth sentence will only appear if the third sentence is also used for the simple reason that the third sentence has a higher score than the fourth sentence. This result is the product of a short document with few sentences. However, in larger documents with more sentences, the cue-phrase list will effectively institute conditions on certain sentence uses. For instance, suppose that the fourth sentence in the above four-sentence paragraph had a higher sentence score than the third sentence. In this case, the fourth sentence is only used if the lower scoring, preceding third sentence is used.

Following the statistical and cue-phrase analysis phases, the document summarizer creates a summary containing the higher ranked sentences which

1 survive the cue-phrase analysis (step 62). The summary may include a conditioned
2 sentence in the event that the relevant condition is satisfied, but will exclude any
3 prohibited sentences. The length of the summary is an author-controlled
4 parameter. From Table 2, a two-sentence summary for the above sample
5 paragraph is as follows:

6
7 Manufacturers have web sites describing
8 their computers. One computer manufacturer
9 offers a money back guarantee.

10
11 The two sentences in the summary had the highest ranking. It is noted that
12 the sentences are organized in the summary according to their order of *appearance*
13 in the document, not in order of their *rank*. In this case, the appearance and rank
14 order are the same, but this does not have to be the case. For example, assume that
15 the third sentence received a higher rank than the second sentence. In the resultant
16 summary, the lower-ranked second sentence would still precede the higher-ranked
17 third sentence because it appears before the third sentence in the document.
18 Ordering a summary based on *rank* reorganizes the author's sentence sequence and
19 might result in a confusing and less readable summary.

20 The two sentence summary did not contain any cue-phrase sentences.
21 However, were the summary expanded to three sentences, it would read as
22 follows:

23
24 Manufacturers have web sites describing
25 their computers. One computer manufacturer

1 offers a money back guarantee. That is why
2 that manufacturer has so many visits to its
3 Internet web site.
4

5 In this summary, the last sentence (i.e., the original fourth sentence) had the
6 third highest sentence score (see Table 2). This sentence also happens to be a
7 conditioned sentence because it contains the phrase "That is why..." which
8 appears on the pre-compiled cue-phrase list. Accordingly, the sentence is used
9 only if the condition is met. In this case, the condition is a depends-on-previous
10 condition, which stipulates that a sentence belonging to this class can be included
11 in a summary only if the preceding sentence is also included. Since the third
12 sentence does appear in the summary, the depends-on-previous condition is met
13 and hence, the fourth sentence can be included in the summary.

14 After the summary is created, the document summarizer displays the
15 summary on the computer monitor in one of four, author-selected UI (user
16 interface) formats (step 64). The first UI format is to insert the summary at the top
17 of the existing document. The document summarizer locates the top of the file,
18 and inserts the summary text before the opening paragraph of the document. Fig.
19 3a shows an existing document 70 with a summary 72 inserted at the top. A
20 second UI format is to create or open a new document and insert the summary in
21 the new document. Fig. 3b illustrates a new document 74 opened and overlaid on
22 an existing document 70. The summary 72 is inserted in the new document 74.

23 The third UI format is to underline or otherwise highlight the important
24 sentences used in the summary. The fourth UI format is to show only the summary
25 sentences without the accompanying text. These third and fourth formats are

1 similar to the conventional presentations described in the Background of the
2 Invention Section.

3 Once the summary is created and displayed to the author, the author can
4 save the summary in the existing document or new document to memory (step 66).

5 A modification of the above computer-implemented method concerns the
6 statistical phase. In the method described above, the content words are counted
7 and all of the sentence scores are derived using the same frequency counts. In
8 some instances, there may be occasions where certain words in the higher ranking
9 sentences unduly dominate and influence the scores of the sentences.

10 The modified technique is an iterative scoring approach. Under this
11 technique, the summarizer initially scores all of the sentences as above on the first
12 iteration. Then, for the next iteration, the summarizer removes the influence of the
13 highest ranking sentence and re-scores the remaining sentences as if the highest
14 ranking sentence was not present. For the next iteration, the influence of the
15 highest scoring sentence found in the previous iteration is removed, and the
16 remaining sentences are again re-scored as if the two highest ranking sentences
17 were not present. This process continues for all of the sentences.

18 To demonstrate this modified statistical analysis, let's apply the analysis to
19 the four-sentence paragraph used above. The first step is to count the content
20 words, while accounting for the stop words and phrase compression. The word
21 count yields Table 1. Next, the sentence scores are derived. The first iteration
22 yields the same score of 2.67 for sentence #2. Here, however, is where the
23 modified method begins to diverge. To remove the influence of the highest
24 ranking sentence, the document summarizer re-computes the sentence scores as if
25 the second sentence were never present in the document. The frequency counts of

the content words are reduced accordingly. Table 4 is a modified version of Table 1 and reflects the absence of the second sentence.

Table 4: Rank of Content Words With Second Sentence Omitted

<u>Content Word</u>	<u>Frequency Count</u>
Computer	3-1=2
Manufacturer	3-1=2
Internet	2
web site	2-1=1
Place	1
Shop	1
Money	1
Visit	1

Next, the remaining three sentences are re-scored using the modified frequency counts for the content words. This results in a ranking of 1.67 for the sentence three, which is second highest.

$$\text{Sentence \#3} = [\text{computer}(2) + \text{manufacturer}(2) + \text{money}(1)] \div 3 \text{ Words} = 1.67$$

The influence of sentence #3 is then removed, and the frequency counts of the content words are reduced accordingly. Table 5 is a modified version of Table 4 and accounts for the absence of the second and third sentences.

**Table 5: Rank of Content Words With
Second and Third Sentences Omitted**

<u>Content Word</u>	<u>Frequency Count</u>
Computer	3-2=1
Manufacturer	3-2=1
Internet	2
web site	2-1=1
Place	1
Shop	1
Money	1-1=0
Visit	1

Continuing this process through the remaining two sentences yields a new sentence rank, given in Table 6.

Table 6: Rank of Sentences With Iterative Re-Scoring Method

<u>Sentences</u>	<u>Score</u>
#2 Manufacturers have web sites describing their...	2.67
#3 One computer manufacturer offers a money back...	1.67
#1 The Internet is a great place to shop for a computer.	1.33
#4 That is why that manufacturer has so many visits to ...	1.00

Notice that using the iterative re-scoring method yields a slightly different sentence ranking with sentence #1 being ranked higher than sentence number #4. A two-sentence summary using the iterative re-scoring method is identical to the

1 two-sentence summary created using the method described above. However, a
2 three-sentence summary is considerably different. A three-sentence summary
3 using Table 6 is as follows:

4
5 The Internet is a great place to shop for
6 a computer. Manufacturers have web sites
7 describing their computers. One computer
8 manufacturer offers a money back guarantee.
9

10 This three-sentence summary is a good example of the situation where the
11 sentences used in the summary are written in order of the *appearance* in the
12 document, and not in order of their *rank*. The beginning sentence in the summary
13 is actually the third highest ranked sentence. Nonetheless, it is written in the
14 summary as the first sentence because it appears in the document before the
15 higher-ranked sentences #2 and #3.

16 In the above example, the counts of the content words appearing in the
17 higher ranking sentences are all reduced by a full count. In other implementations,
18 the frequency counts can be changed by varying degrees depending upon the
19 degree of influence introduced by the higher ranking sentences the manufacturer or
20 author desires to remove. For instance, the summarizer might compensate by
21 subtracting a fractional amount (say, 0.3 or 0.5) from each count corresponding to
22 words that appear in the highest ranking sentence. Alternatively, the compensation
23 amount might vary depending upon whether the content word has a high or low
24 frequency count compared to other content words. The amount that word counts
25 are compensated during this dynamic scoring process can be determined and set by

1 the manufacturer or author according to various statistical or mathematical
2 approaches which appropriately negate the influence of the content words
3 appearing in the higher ranking sentences.

4 The document summarizer is advantageous over prior art summarizers
5 because it is designed from the author's standpoint. It enables authors to
6 automatically create summaries of their writings using a combined statistical and
7 cue-phrase approach. Once created, the summarizer presents a UI that enables the
8 author to place the summary at the top of the document or in a new document.
9 This placement is convenient and useful to the author. The author is then free to
10 revise the summary as he/she wishes.

11 Another advantage of the document summarizer stems from the combined
12 statistical and cue phrase processing. This dual analysis is beneficial because the
13 statistical component ensures that a summary will always be produced, and the cue
14 phrase component improves the quality of the resulting summary.

15 In compliance with the statute, the invention has been described in language
16 more or less specific as to structure and method features. It is to be understood,
17 however, that the invention is not limited to the specific features described, since
18 the means herein disclosed comprise exemplary forms of putting the invention into
19 effect. The invention is, therefore, claimed in any of its forms or modifications
20 within the proper scope of the appended claims appropriately interpreted in
21 accordance with the doctrine of equivalents and other applicable judicial doctrines.